

基于代价敏感和近似分类质量的决策粗糙集属性约简研究<sup>\*</sup>陈婉清, 秦亮曦<sup>†</sup>

(广西大学 计算机与电子信息学院, 南宁 530004)

**摘要:**针对决策粗糙集属性约简在引入代价后分类精度不高的问题,对其中代价敏感与分类精度的平衡进行了研究。将分类总代价和近似分类质量作为属性约简过程中的约束条件,结合模拟退火方法,提出了一个基于代价敏感和近似分类质量的决策粗糙集属性约简(ARACOQ)算法。利用 UCI 数据集对算法进行了模拟实验,实验结果验证了 ARACOQ 算法的有效性,该算法能够在可承受代价范围内找到一个分类精度最高的属性约简集。

**关键词:** 决策粗糙集; 属性约简; 代价敏感; 近似分类质量; 分类精度

**中图分类号:** TP391

## Study on DTRS attribute reduction constrained by cost-sensitive and classification quality

Chen Wanqing, Qin Liangxi<sup>†</sup>

(School of Computer, Electronics &amp; Information Guangxi University, Nanning 530004, China)

**Abstract:** Aiming at the low precision problem while the cost is introduced into attribute reduction of decision-theoretic rough set, it is studied the balance between the total cost and the precision in classification. The total cost of the classification and the approximate classification quality are used as the constrained criteria in the attribute reduction procedure, combined with simulated annealing method, it is proposed a DTRS attribute reduction algorithm constrained by cost-sensitive and classification quality (hereinafter referred as ARACOQ). The simulation experiments are carried out by using UCI data set, the results verify the effectiveness of the ARACOQ algorithm, it can find an attribute reduction set with the highest classification precision within the affordable cost range.

**Key Words:** decision-theoretic rough set; attribute reduction; cost sensitive; classification quality; precision

## 0 引言

粗糙集模型<sup>[1]</sup>是由 Pawlak 于 1982 年提出的一种计算工具<sup>[2]</sup>,主要用于分析和处理不精确性和模糊性的数据。经典粗糙集理论是基于严格的代数包含关系建立的,但是在实际应用中往往难以满足精确的代数包含,因此导致经典粗糙集在处理实际分类问题时缺乏容错能力。基于此,Yao 等人将贝叶斯风险相关理论与粗糙集相结合,提出了具有容错能力的决策粗糙集模型<sup>[3]</sup>。

随着决策粗糙集研究的逐步深入,其属性约简问题得到了学者的广泛关注。属性约简的目的是在保证信息系统某些关键特征值不变的情况下,将冗余的属性删减。决策粗糙集中决策语义的引入,导致其属性约简不再具有单调性。因此,Yao 等人首先提出了决策粗糙集的属性约简问题,并提出了一种基于属性  $\alpha$ -正域重要度的属性约简<sup>[4]</sup>方法;贾修一等人将决策风险最小化作为最优化目标,提出了一种基于决策风险最小化的属性约简<sup>[5]</sup>方法;Bi 等人从代数理论和信息论两个方面提出了基

于最小决策代价的属性约简方法<sup>[18]</sup>;Zhang 等人在决策粗糙集模型下,提出了不完备系统的最小代价属性约简方法<sup>[19]</sup>;Song 等人将决策粗糙集与模糊集结合,提出了两种属性约简方法:全局约简和局部约简,全局约简能保持所有决策类的代价不变或减少,而局部约简则能保持单个决策类的代价不变或减少<sup>[20]</sup>。

通常情况下,不同的测试属性集会带来不一样的分类结果。在一定的范围内,测试属性集中属性个数越多,错误分类的结果越少,则误分类代价越小,分类精度越高<sup>[6]</sup>。然而在日常生活中,数据的获取需要花费一定的经济或时间成本,即测试代价<sup>[7]</sup>。例如,在医疗诊断中,各种医疗检测都需要花费一定的费用<sup>[8]</sup>。随着测试属性个数的增多,在误分类代价减少的同时,也会使测试代价增加。因此,在实际问题中,需将测试代价和误分类代价同时考虑,并找到兼顾两者的一个平衡点。基于此,Min 等<sup>[7]</sup>在粗糙集属性约简问题中,率先引入测试代价作为约束条件。李华雄等人将代价敏感引入决策粗糙集,提出了代价敏感的决策风险最小化属性约简方法<sup>[6]</sup>。刘偲等人采用模拟退火算法结合传统决策粗糙集正域约简算法,搜索测试代价总和最小

**基金项目:** 国家自然科学基金资助项目(61363027);广西重研发点计划项目桂科(AB16380260)

**作者简介:** 陈婉清(1993-),女,广西南宁人,硕士研究生,主要研究方向为决策粗糙集、数据挖掘、机器学习;秦亮曦(1963-),男(通信作者),广西灵川人,教授,博士,主要研究方向为数据挖掘、决策粗糙集、深度学习等(qin\_lx@126.com)。

的正域约简属性集<sup>[9]</sup>, 取得了较好的结果。

在以上研究中, 一些没有考虑代价, 得到的是满足某些条件(如基于属性  $\alpha$ -正域重要度的属性约简<sup>[4]</sup>)的属性子集; 而一些考虑了代价, 包括误分类代价、测试代价或包含两者的总代价, 得到的是具有最小代价的属性子集, 但此类属性集的分类精度往往不高。在实际问题中, 分类精度应该是首先要考虑的问题, 如在诊断重大疾病时, 诊断精确性的地位是远远高于测试代价的。因此, 在决策粗糙集属性约简的问题上, 分类精度和分类代价应该综合考虑, 应在可承受的分类总代价范围内, 尽可能提高分类精度。

在粗糙集理论中, 近似分类质量表明了应用知识  $R$  能确切地划入已知分类的对象的百分比<sup>[10]</sup>。近年来, 基于近似分类质量的属性约简算法也不断提出<sup>[14-17]</sup>。在近似分类质量不变的前提下求约简集, 可以保证约简集的分类决策能力不会被减弱<sup>[14]</sup>。因此, 本文将近似质量作为属性约简的约束条件之一, 以此保证分类决策能力不会被大幅度的削减。

本文在属性约简问题中将分类总代价和近似分类质量作为属性约简的迭代准则, 在可承受的分类总代价的范围内, 寻找对分类决策能力高的测试属性子集。

## 1 决策粗糙集的基本概念

设  $\Omega = \{w_1, w_2, w_3, \dots, w_s\}$  表示  $s$  个状态的集合;  $A = \{a_1, a_2, a_3, \dots, a_m\}$  表示  $m$  个可能的决策;  $x$  表示为论域中的某对象;  $\tilde{x}$  表示为对象  $x$  的属性特征描述;  $P(w_j | \tilde{x})$  表示在  $\tilde{x}$  描述下的对象  $x$  具有  $w_j$  状态的条件概率;  $\lambda(a_i | w_j)$  表示在  $w_j$  状态下的作出  $a_i$  决策的风险代价, 其中  $\lambda$  通常是由经验得出。在决策粗糙集模型中, 通常考虑 3 种分类决策, 即正域(属于  $w_j$  类)、负域(不属于  $w_j$  类)和边界域(待定结果)。因此, 对于具有  $\tilde{x}$  描述的对象  $x$ , 在采取  $a_i$  决策情况下, 可能带来的决策风险期望为

$$R(a_i | \tilde{x}) = \sum_{j=1}^s \lambda(a_i | w_j) P(w_j | \tilde{x})$$

为了叙述方便, 本文只考虑包含两种互补的状态的集合  $\Omega = \{X, \sim X\}$ 。设决策集  $A = \{a_P, a_N, a_B\}$ , 其中  $a_P$ ,  $a_N$ ,  $a_B$  分别表示决策为正域  $POS(X)$ 、负域  $NEG(X)$  和边界域  $BND(X)$ 。当对象  $x \in X$  时, 对象划分到相应区域  $POS(X)$ 、 $NEG(X)$  和  $BND(X)$  的损失函数为  $\lambda_{PP}$ 、 $\lambda_{NP}$  和  $\lambda_{BP}$ 。反之, 当对象  $x \notin X$  时, 对象划分到相应区域  $POS(X)$ 、 $NEG(X)$  和  $BND(X)$  的损失函数为  $\lambda_{PN}$ 、 $\lambda_{NN}$  和  $\lambda_{BN}$ 。在粗糙集理论中, 等价类  $[x]_R$  表示具有完全相同特征描述的对象  $x$ 。因此, 本文将具有特征  $\tilde{x}$  描述的对象  $x$  用等价类  $[x]_R$  来表示。因此, 可以得到 3 种决策的期望风险为

$$R(a_P | [x]_R) = \lambda_{PP} P(X | [x]_R) + \lambda_{PN} P(\sim X | [x]_R)$$

$$R(a_N | [x]_R) = \lambda_{NP} P(X | [x]_R) + \lambda_{NN} P(\sim X | [x]_R)$$

$$R(a_B | [x]_R) = \lambda_{BP} P(X | [x]_R) + \lambda_{BN} P(\sim X | [x]_R)$$

对于决策代价函数值的大小, 显然有如下的关系:

$$\lambda_{PP} \leq \lambda_{BP} < \lambda_{NP}, \quad \lambda_{NN} \leq \lambda_{BN} < \lambda_{PN}$$

考虑到正确分类的风险为 0, 即  $\lambda_{PP} = \lambda_{NN} = 0$ , 因此, 上述的三种分类的期望风险可以表示为

$$R(a_P | [x]_R) = \lambda_{PN} P(\sim X | [x]_R)$$

$$R(a_N | [x]_R) = \lambda_{NP} P(X | [x]_R)$$

$$R(a_B | [x]_R) = \lambda_{BP} P(X | [x]_R) + \lambda_{BN} P(\sim X | [x]_R)$$

根据贝叶斯最小风险决策原则, 可以得到决策规则如下:

如果  $R(a_P | [x]_R) \leq R(a_N | [x]_R)$  并且  $R(a_P | [x]_R) \leq R(a_B | [x]_R)$ , 那么选择  $x \in POS(X)$ 。

如果  $R(a_N | [x]_R) \leq R(a_P | [x]_R)$  并且  $R(a_N | [x]_R) \leq R(a_B | [x]_R)$ , 那么选择  $x \in NEG(X)$ 。

如果  $R(a_B | [x]_R) \leq R(a_P | [x]_R)$  并且  $R(a_B | [x]_R) \leq R(a_N | [x]_R)$ , 那么选择  $x \in BND(X)$ 。

## 2 模拟退火算法概述

模拟退火算法是 Metropolis 等人于 1953 年提出的一种随机优化算法。该算法通过模拟热力学中物体从高温开始, 缓慢地降温(这个过程被称为退火), 最终在某一温度达到热平衡的过程, 从而求解优化问题的最小值。已经证明, 虽然温度下降缓慢, 但模拟退火算法最终一定能达到全局最优<sup>[12, 13]</sup>。

从设定的初始温度  $T_0$  及初始状态  $x(0)$  开始, 模拟退火算法随机地从可行解中, 持续进行“产生新解-判断-接受/舍弃”的迭代过程, 从而产生一个状态序列  $x(0), x(1), \dots, x(i)$ , 且新状态  $x(i+1)$  只依赖于前一个状态  $x(i)$ , 与前面的状态  $x(0), x(1), \dots, x(i-1)$  无关, 因此该状态序列构成一个马尔可夫链。而模拟退火算法实际上是通过马尔可夫链的演化过程, 逐步逼近问题的最优解。到达停止准则后, 使用衰减函数减少控制参数的值, 重复以上步骤, 当控制参数到达终止时, 即得到最优解。

## 3 代价敏感与属性约简

### 3.1 测试代价

假设各个样本中的同一属性值的测试代价相同, 则在测试属性集  $B$  上计算的样本  $x$  的测试代价等于  $B$  中每个属性  $c_i \in B$  测试代价的总和, 测试代价设为一个非负实数。因此, 可得计算测试代价函数<sup>[6]</sup>:

$$Test\ cost(x, B) = \sum_{i=1}^{|B|} TC(c_i)$$

其中:  $TC(c_i)$  为  $B$  中单个属性集的测试代价。

### 3.2 误分类代价

设决策表信息系统  $S = (U, At = B \cup D, V, f)$ , 其中  $U = \{x_1, x_2, \dots, x_k\}$  为非空子集论域,  $B$  和  $D$  分别为条件属性集和决策属性集, 给定条件属性子集  $B \subseteq B$ , 由  $B$  确定的等价关系  $R_B$  和样本  $x \in U$  在属性集  $B$  上的等价类为<sup>[6]</sup>:

$$R_B = \{(x, y) \in U \times U \mid \forall a \in B, I_a(x) = I_a(y)\}$$

$$[x]_B = \{y \in U \mid (x, y) \in R_B\}$$

根据决策粗糙集的决策规则可以计算出不同属性子集的误分类代价值, 计算形式如下:

当对象  $x$  被划分到正域时,

$$Errorcost(x, B) = \lambda_{PN} P(\sim X | [x]_B)$$

当对象  $x$  被划分到负域时,

$$Errorcost(x, B) = \lambda_{NP} P(X | [x]_B)$$

当对象  $x$  被划分到边界域时,

$$Errorcost(x, B) = \lambda_{BP} P(X | [x]_B) + \lambda_{BN} P(\sim X | [x]_B)$$

### 3.3 代价敏感

设  $Sumcost(x, B)$  为样本  $x$  在测试属性集  $B$  上的分类总代价, 分类总代价就是误差代价与测试代价之和, 即

$$Sumcost(x, B) = Errorcost(x, B) + Testcost(x, B)$$

### 3.4 近似分类质量

设  $\gamma(x, B)$  为样本  $x$  在测试属性集  $B$  上的近似分类质量<sup>[10]</sup>, 则有

$$\gamma(x, B) = \frac{|POS_{[x]_B}(D)|}{|U|}$$

### 3.5 属性约简

文献[6]中采用启发式算法, 每次均将具有最小的分类总代价的属性加入最优属性集中。通常情况下, 同一样本下获得的最优属性集是一致的。然而分类总代价最小的属性, 分类精度可能不高。因此, 文献[6]中给出的算法是有局限性的。本文将模拟退火方法<sup>[11]</sup>引入决策粗糙集的属性约简问题, 综合考虑分类总代价和分类精度之间的关系, 即在可承受的分类总代价的范围内(本文将全属性集分类总代价的 10% 作为最大可承受的分类总代价), 用近似分类质量最大限度的保证分类决策能力的提高。其中, 最大可承受分类代价为

$$AffordSC = Sumcost(x, B_{Oringe}) \times 10\%$$

使用模拟退火算法随机找到一组属性子集, 并判断该属性子集是否满足以下条件, 即

$$0 < Sumcost(x, B) \leq Max(AffordSC) \\ \& \& isMax(\gamma(x, B))$$

其中:  $x$  为样本对象,  $Sumcost(x, B)$  为在属性子集  $B$  下对象  $x$  的分类总代价,  $isMax(\gamma(x, B))$  判断属性子集  $B$  的近似分类质量是否为已找到的属性子集中最高的。

当满足上述条件, 则该属性子集为最优属性集。由此给出如下算法。

算法 1 基于代价敏感的决策粗糙集属性约简算法(以下简称 ARACQ 算法)

输入: 一个决策表  $S=(U, At=B_{Oringe} \cup D, V, f)$ , 待分类样本  $x$ , 样本的全属性集  $B_{Oringe}$ , 误差代价矩阵, 测试代价矩阵, 模拟退火算法的参数设置如下: 马尔可夫链的迭代次数  $K=1000$ , 马尔可夫链长度  $MarkovLength=1000$ , 衰减因子  $DecayScale=0.95$ , 步长  $StepFactor=0.02$ , 初始温度  $t=30$ , 容差  $t_{min}=10^{-8}$ 。

输出: 最优属性子集  $BestB$ , 分类总代价  $BestSC$ , 近似分类质量  $BestQuality$  (初始值为 0)。

a) 计算全属性分类总代价  $Sumcost(x, B_{Oringe})$  和可承受分

类总代价  $AffordSC$ 。

b) 初始化  $t_{min}$ ,  $t$ ,  $B=B_{Oringe}$ ,  $Num=0$ 。

c) 通过随机添加、替换、删除属性的方式, 由属性集  $B$  产生新的属性子集  $B^+$ , 那么  $B=B^+$ ,  $Num=Num+1$ 。

d) 计算属性子集  $B$  的分类总代价和近似分类质量。若  $Sumcost(x, B) < AffordSC \& \& \gamma(x, B) > BestQuality$ , 那么  $BestB=B$ ,  $BestSC=Sumcost(x, B)$ ,  $BestQuality=\gamma(x, B)$ 。

e) 判断当  $t > t_{min}$  并且结果不收敛时, 返回 c), 同时  $t=DecayScale \times t$ 。

f) 若  $Num < 5$ , 返回 c)。

g) 输出  $BestB$  即为分类总代价在可承受范围内分类精确度最高的最优测试属性集。

ARACQ 算法结束条件有两个, 其一是初始温度  $t$  衰减到规定的值, 其二是若在马尔可夫链  $K$  次迭代过程中, 实验结果没有发生任何改变, 即结果收敛。若结果一直未收敛, 当温度  $t$  也会衰减到规定的值时结束该算法, 并且已有相关论文证明模拟退火算法是可收敛的, 只不过收敛速度较慢<sup>[12]</sup>。

在运算时间上, 启发式算法较 ARACQ 算法有明显优势, 但是在分类精确度上, ARACQ 算法较启发式算法具有绝对优势。并且在实际问题中, 在一定的分类代价范围内, 分类精度的提高是非常重要的。

## 4 实验结果及分析

为了验证 ARACQ 算法的有效性, 本文使用 UCI 的数据集对算法进行了模拟实验, 并与李华雄提出的决策粗糙集代价敏感属性选择及分类算法进行了对比分析。

模拟退火算法属于随机算法, 每次运行得到的结果可能不一致, 因此在每个数据集上, 本文均做 10 次实验, 取得的平均值作为运行结果。

实验的机器为 Intel<sup>(R)</sup> Xeon<sup>(R)</sup> 的 3.50 GHz CPU、内存为 8 GB, 64 位的 Windows10 操作系统, 算法在 MATLAB 平台上实现。本实验使用的数据集均来自 UCI 数据库, 这三组数据集分别为 Car、WPBC (breast-cancer-Wisconsin (diagnostic)) 和 Spambase。针对数据集中少量数据缺失的情况, 使用最频值填充法进行补齐。删除 WPBC 数据集中的 ID 列, 并对数据进行归一化和离散化。在数据集 Car 中, 将类别数调整为 2 个, 即将“good”和“vgood”均归于“acc”类别中。实验数据进行预处理后的基本信息如表 1 所示。

表 1 实验数据基本信息

名称	类别数	属性个数	处理后属性数	样本数
Car	2	6	6	1728
WPBC	2	34	33	198
Spambase	2	57	57	4601

在实验中, 假定误分类代价满足  $\lambda_{PP} \leq \lambda_{BP} < \lambda_{NP}$  和  $\lambda_{NN} \leq \lambda_{BN} < \lambda_{PN}$ ,  $\lambda_{PP} = \lambda_{NN} = 0$ 。为了实验具有对比性, 对第一组数据集 (Car) 和



后两组数据集 (WPBC 和 Spambase) 采用不同的误分类代价矩阵, 如表 2 和 3 所示。

表 2 Car 的误分类代价矩阵

名称	POS	BND	NEG
X	0	2	8
$\sim X$	9	2	0

表 3 WPBC 和 Spambase 的误分类代价矩阵

名称	POS	BND	NEG
X	0	8	20
$\sim X$	15	7	0

在 UCI 数据中, 并未给出测试代价。因此, 本文假定每个数据集中的属性测试代价服从正态分布  $N(\mu, \sigma)$ 。设 Car 和 Spambase 属性的测试代价服从  $N(0.2, 0.1)$ , WPBC 属性的测试代价服从  $N(0.02, 0.01)$ , 在 Matlab 中使用正态随机函数为各个数据集的属性生成测试代价。

由于本文将分类近似质量作为其中的一个指标, 因此在选取训练样本时, 需要截取该样本数据每一个决策类别均覆盖的数据作为训练样本。在 car 集上, 选取第 1028 至 1227 条数据作为训练样本。在 Spambase 集上, 选取第 1714 至 1913 条数据作为训练样本。在 WPBC 数据集上, 选取前 100 条数据作为训练样本。其余的样本均为测试样本<sup>[6]</sup>。在三个数据集上, 分别使用 ARACOQ 算法和李华雄的启发式算法计算得到最优测试属性集, 然后使用 WEKA 在测试样本上验证, 最优测试属性集

的分类精度。本文的实验结果将通过代价约简率、分类精度和属性约简率来做比较<sup>[6]</sup>。其中:

$$\text{代价约简率} = \frac{\text{全属性总代价} - \text{最优属性总代价}}{\text{全属性总代价}}$$

$$\text{属性约简率} = \frac{\text{全属性数} - \text{最优属性集大小}}{\text{全属性数}}$$

从表 4 中可以发现, 两种算法的代价约简率均达到了 90% 以上。ARACOQ 算法和启发式的代价约简率在同一个数据集上相差不会超过 7%。因此, 证明了 ARACOQ 算法能够大大降低分类总代价总和。在实际问题中, 能够很好地解决因代价过高而耽误病情等问题。

表 5 展示了全属性集、启发式算法得到的最优属性集和 ARACOQ 算法得到的最优属性集, 使用 WEKA 对测试样本集进行风险最小化的决策粗糙集决策分类得到的分类平均精度。

综合分析表 4 和 5 的数据结果可以发现, 一味地追求总代价最小化, 会导致分类精度大幅度地下降, 因此需要将两者结合起来考虑。与启发式算法相比, 在同一属性集上, 当分类总代价不超过最小代价的 7% 时, 分类精度得到了大幅度的提高。尤其在 Car 和 Spambase 这些测试样本数量很多的数据集上, 效果更为显著。证明了 ARACOQ 算法能够在可承受的分类总代价的范围内, 能够大幅度地提高了分类精度。与全属性集相比, ARACOQ 算法在代价约简率均达到 90% 的同时, 分类精度也维持在了 10% 以内的下降范围。在实际问题中需要综合考虑代价和分类精度之间的关系, 因此 ARACOQ 算法的思想和实验结果更加符合实际。

表 4 平均总分类代价

数据集	全属性总代价	启发式算法总代价 (平均)	ARACOQ 算法总代价 (平均)	启发式算法代价约简率 (平均) / %	ARACOQ 算法代价约简率 (平均) / %
Car	1904.10	2.45	33.17	99.87	98.26
WPBC	66.54	2.06	4.36	96.90	93.45
Spambase	23648.00	6.22	1567.33	99.97	93.37

表 5 分类的平均精度

属性集	全属性分类精度 / %	启发式算法最优属性集分类精度 / %	ARACOQ 算法最优属性集分类精度 (平均) / %
Car	93.20	50.70	84.37
WPBC	72.00	66.00	70.67
Spambase	84.80	67.30	76.64

表 6 属性平均约简率

数据集	处理后全属性集	启发式算法约简属性集 (平均属性个数)	ARACOQ 算法约简属性集 (平均属性个数)	启发式算法约简率 (平均) (%)	ARACOQ 算法约简率 (平均) (%)
Car	6	1	2.7	83.33	55.00
WPBC	33	3	2.6	90.91	92.12
Spambase	57	1	19.7	98.25	65.44

由表 6 可以看出, 在全属性集个数较少的情况下, ARACOQ

算法和启发式算法的约简率相差较多。在全属性集数量个数适

中的情况下, ARACOQ 算法在大幅度提高分类精度的同时, 约简率比启发式算法高。在全属性集个数较多的情况下, ARACOQ 算法与全属性集相比约简率达到了 65.44%。证明了 ARACOQ 算法具有很强的属性约简能力, 能够大幅度减少属性个数, 获得最优测试属性集。

图 1~3 分别用图形的方式展示了归一化的平均总测试代价 (由于数量级不同, 将三组数据使用同一个数据量展示)、分类的平均精度和归一化属性集个数。

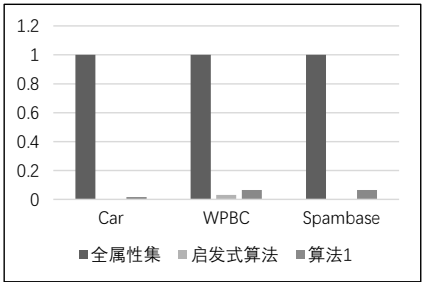


图 1 归一化平均总分类代价

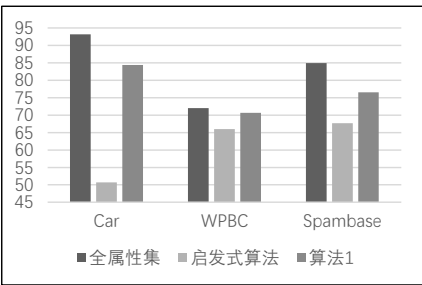


图 2 分类的平均精度

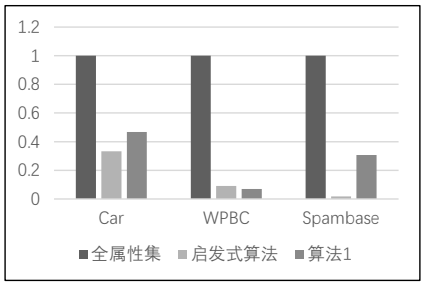


图 3 归一化属性集个数

综合分析图 1 和图 2 的图表信息可以更清楚地看出, 在同一数据集上, ARACOQ 算法比启发式算法具有更高的分类精度, 并且分类总代价相差不大。ARACOQ 算法的分类总代价控制在可承受的分类代价范围内, 这个是通过人工设定的, 可以设置为总分类代价的 10% 等。因此, 在实际问题中, 具有比启发式算法更强的适用性。

从图 2 和 3 中可以直观地观察到, ARACOQ 算法得出的最优属性集与全属性集相比, 属性个数大幅度地下降, 大大降低了分类总代价总和并且分类精度相差不大。

综上所述, ARACOQ 算法能够在一定承受代价范围内, 给出分类精度较高的最优测试属性集。与全属性集相比,

ARACOQ 算法得到的最优属性集属性个数和分类总代价均大幅度地减少。与启发式算法相比, ARACOQ 算法得到的最优属性集分类正确度大幅度提高, 具有非常重大的现实意义。

## 5 结束语

在决策粗糙集属性约简中, 传统方法以近似分类质量作为衡量的标准, 没有考虑分类代价, 从而可能导致分类代价过高; 而一些算法以分类代价作为优化的目标, 追求代价最小化从而会导致分类精度不高。

因此, 本文综合考虑分类精度和分类总代价之间的平衡, 提出了一种基于分类总代价和近似分类质量的决策粗糙集属性约简方法。该方法通过模拟退火算法随机找到一组分类总代价在可承受范围内的属性集, 再通过比较近似分类质量, 找到最优属性集, 使得分类结果兼具分类总代价适中且分类精度较高的特性。实验结果表明本文提出的算法是有效的。

在要求高精度的实际问题, 如临床诊断中, 使用 ARACOQ 算法能够有效的删减检查项目数, 节约诊断成本, 并且确保较高精度的诊断结果。相较于启发式算法, ARACOQ 算法时间复杂较高, 所以运行时间方面明显弱于启发式算法。并且 ARACOQ 算法对于不同的初始值, 得到的最优属性集会略有不同, 分类总代价和分类精度也会相应发生一些变化。因此, 不同初始值在同一数据集上的变化规律和如何确定不同数据集的初始值将成为下一步的讨论重点。

## 参考文献:

- [1] Pawlak Z. Rough Sets [J]. International Journal of Computer and Information Sciences, 1982, 11 (5), 341-356.
- [2] 王国胤, 姚一豫, 于洪. 粗糙集理论与应用研究综述 [J]. 计算机学报, 2009, 32 (7), 1229-1246.
- [3] Yao Yiyu. Decision-theoretic rough set models [M]// Yao J, Lingras P, Wu W Z, et al. Rough set and Knowledge Technology. Heidelberg: Springer, 2007: 1-12.
- [4] YaoYiyu, ZhaoYan. Attribute reduction in decision-theoretic rough set models [J]. Information Sciences, 2008, 178 (17): 3356-3373.
- [5] 贾修一, 商琳, 陈家骏. 决策风险最小化属性约简 [J]. 计算机科学与探索, 2011, 5 (2): 155-166.
- [6] 李华雄, 周献中, 黄兵, 等. 决策粗糙集与代价敏感分类 [J]. 计算机科学与探索, 2013, 7 (2): 126-135.
- [7] Liao Shujiao, Zhu Qingxin, Min Fan. Cost-sensitive attribute reduction in decision-theoretic rough set models [J]. Mathematical Problems in Engineering, 2014 (2): 1-9.
- [8] Ju H R, Yang X B, Yu H L, et al. Cost-sensitive rough set approach [J]. Information Sciences, 2016, 355 (C): 282-298.
- [9] 刘偲, 秦亮曦. 测试代价敏感的决策粗糙集正域约简 [J]. 计算机科学与探索, 2017 (6): 1014-1020.
- [10] 朱晓钟, 杨勇, 朱英丽. 一般关系粗糙集的近似分类精度和质量 [J].

计算机应用与软件, 2011, 28 (5): 52-54.

[11] 贾修一, 商琳. 一种求三支决策阈值的模拟退火算法 [J]. 小型微型计算机系统, 2013, 34 (11): 2603-2606.

[12] 姚新, 陈国良. 模拟退火算法及其应用 [J]. 计算机研究与发展, 1990 (7): 1-6.

[13] 卢宇婷, 林禹攸, 彭乔姿, 等. 模拟退火算法改进综述及参数探究 [J]. 大学数学, 2015, 31 (6): 96-103.

[14] 吴陈, 李丹丹. 基于粗糙集的关联规则挖掘方法的研究与应用 [J]. 电子测量技术, 2016, 39 (7): 44-48.

[15] 高晓红. 基于分类质量的决策系统属性约简新算法 [J]. 计算机与现代化, 2010 (1): 19-22.

[16] 杨成福, 舒兰. 基于近似分类质量的决策表属性约简算法 [J]. 河西学院学报, 2006, 22 (5): 1-3.

[17] 徐德友, 胡寿松. 一种基于粗糙集的近似质量求取属性约简的决策算法 [J]. 控制与决策, 2003, 18 (3): 313-316.

[18] Bi Zhongqin, Xu Feifei, Le Jingsheng, et al. Attribute reduction in decision-theoretic rough set model based on minimum decision cost [J]. Concurrency & Computation Practice & Experience, 2016, 28 (15): 4125-4143.

[19] Zhang Yimeng, Jia Xiuyi, Tang Zhenming. Minimum cost attribute reduction in incomplete systems under decision-theoretic rough set model [C]// Proc. of the 12th International Conference on Natural Computation, Fuzzy Systems and Knowledge Discovery. 2016: 940-944.

[20] Song Jingjing, Tsang E C C, Chen Degang, et al. Minimal decision cost reduct in fuzzy decision-theoretic rough set model [J]. Knowledge-Based Systems, 2017, 126 (C): 104-112.